# Medical literature's semantic status and solution

## Zhang Jianping and Liu Yao

Institute of Scientific and Technical Information of China

No.15, Fuxing Road, Haidian District, Beijing, China

zjping51@126.com

liuy@istic.ac.cn

ABSTRACT. *This article presents the status of literature's semantic, especially the medical literature, and the insufficient of semantic annotation tools. Facing with those problems, we pretend to use NLP techniques and machine learning methods to construct skin disease ontology automatically to solve these problems.*

**Keywords:** Semantic Annotation, Medical Literature, Natural Language Processing

1. **Introduction.** Web founder Tim Berners-Lee put forward the concept of the Semantic Web in 1998 , with the W3C officially launched Semantic Web Activity in February 2001, semantic retrieval became the mainstream of network research and development. Facing of the geometric published biomedical literature (BIO), how to reveal the semantic information from the text data is becoming important. Semantic, is to choose the appropriate semantic tags, with what semantic features reflected, so that information becomes a machine-processable information. There are two key issues involved in Semantics: the generation of semantic metadata and semantic annotation.

Semantic metadata (also known as Tag Ontology) provides the data's semantic information, it is also an important tool for organizing library semantic resources. Semantic metadata's generation technology and domain knowledge base or ontology's construction, from the technology point of view, they are consistent. Its technology is quite mature. Hanxian Pei[1] et al. put forward a Wikipedia-based construction method, the target semantic metadata was extracted by analyzing the table of contents entries in a category (table-of-contents), and according to the title of the section of the document to give the instance a semantic markup correctly, and ultimately build a training corpus. Solve two problems in building automatic semantic metadata's generation system: target semantic metadata extraction and corpus construction, laid the foundation for further large-scale applications of semantic metadata. Liu Yao et al[2] proposed the use of natural language processing (NLP) theory and technology has been generally accepted in the field of knowledge, such as professional thesaurus, professional dictionary, professional teaching or

authoritative book reconstruction and use; with the help of experts in the field of knowledge, web-based knowledge acquisition and processing; establish a limited text ontology self-learning mechanism, in order to achieve field ontology concept describe the system automatically constructed.

Semantic annotation is marking text or symbols of the original data, to have semantic information, not only the man can understand, but also the machine can understand. Semantic Web is generally in XML markup language data for label , described as a data model in RDF / XML, and combining Ontology, the data is marked with a clear meaning, so that the machine can understand. Semantic annotating tool today has been formed there such as Gate ， OntoMat Annotizer, Seeker, MnM SHOE Knowledge Annotator, Annotea SMORE SemanticWord, Melita, SemanticMarkup Plug-In Briefing Associate Yawas Annozilla 10 [3], they are all ontology-based annotating tool. Now, there are normaly three Semantic annotation methods[4] : (1) Artificial annotation, by experts to determine the applicable set of concepts, metadata elements, the establishment of the semantic data in RDF or HTML language tags. This process requires a lot of manpower and material resources. (2) the use of Document Type Definition (DTD) and document mode (Schema) concept mapping and annotating, because of the document's structure and element was definited detailed by the SGML/XML document, so that after the establishment of specific concepts' special mapping of the DTD/Schema, the SGML/XML document's DTD/Schema's content element tag can be transformed to homologous concept element tag automaticly, but this method also needs to be checked and modified manually; (3)annotation with analysis of the vocabulary's meaning ,the method mainly uses automaticly analysis and extracting the vocabulary to establish the mapping of the vocabulary and the concept's classification, then find the document's concept classification, even its annotation relationship with other class, to annotate with these concepts. But this method has domain restriction and also needs to be checked manually. Now, there are more mature Semantic annotating Platform such as AeroDAML, MnM, MUSE, Armadillo, KIM, SemTag, Ont-O-Mat, etc.

2. **Problem Statement and Preliminaries.** To reducing the man's participate in the generation of semantic metedata, some semi-automatic assistance skills were used to these annotating tools. In the part of constructing system, CREAM system's expanding researches of S_CREAM[5] and PASNKOW[6], S_CREAM achieved semi-automaticly annotation of the web, which is based on Amilcare's information extraction. PANKOW （Pattern-based Annotation through Knowledge on the Web） uses  a studying method based on unsupervised framework to classify the example to ontology, at first processes the web page, find nouns may be became ontology ,then accord to a schema to generate phrases contained in the noun, using Google API to find the times of the phrase, use the calculated weight to determine the noun's classification in the ontology system. To be annotated automatically, a lot of work focus on designing semantic annotation metadata 's generate automatically model. Hung[7]et al. Based on Web resource mining, acquires training corpora necessary to describe both the thematic categories and the metadata

extracted from the texts. The approach then finds the corresponding relationships among them by means of categorization and thus generates thematic metadata for the textual data. Yang and Lee[8]proposed a method using machine-learning to generate web page's semantic metedata automatically. The proposed automated process adopts the self-organizing map algorithm to cluster training Web pages and conducts a text mining process to discover some semantic descriptions about the Web pages. A.Dingli[9]proposed a framework of Armadillo. The methodology is based on a combination of in-formation extraction, information integration and machine learning techniques. Learning is seeded by extracting information from structured sources (e.g. databases and digital libraries). Retrieved information is then used to partially annotate documents. These annotated documents are used to bootstrap learning for simple Information Extraction (IE) methodologies, which in turn will produce more annotations used to annotate more documents. H.Graubitz[10] displayed DIAsDEM framework , it is based on cluster text units into homogeneous groups whose labels form the XML tags surrounding their contents and from which the document type definition for the collection is derived. J.Li[11]described a mechine-learing method, it's semantic annotation is based on sentences, and converted to RDF, it adopt a dependency grammar – Link Grammar for this purpose Dill[12]describes Seeker, a platform for large-scale text analytics, and SemTag, an application written on the platform to perform automated semantic tagging of large corpora. We apply SemTag to a collection of approximately 264 million web pages, and generate approximately 434 million automatically disambiguated semantic tags, published to the web as a label bureau providing metadata regarding the 434 million annotations.

Medical literature's semantic metadata's generation is similar of metadata. Liu yao[13]proposed a construction method of Chinese medicine ontology based on history literature, use of history literature and authority domain knowledge to illustrate Chinese medicine's core concept. To construct and obtain Chinese medical knowledge system automatically. Meeraman ,R[14]proposed to construct ontology based on NLP, this method uses WWW to analysis domain literature, can ignore linguistic knowledge and resource ,suggest a linguistic friendly interface to enhance language processing based on ontology such as semantic annotation. Deng ziping[15]proposed some methods to construct domain ontology manually. With the source of electronic medical record for research, focus on solving three key question of construction of medical domain ontology, research on medical diagnosis and treat ontology system's automatically generate theory and technology. Fang an[16]and Fu qiang[17]' research is based on Hadzic of Australian who proposed four-dimensional medical ontology model , improve it to construct HFMD (Hand, Foot and Mouth Disease, )and arrhythmia ontology. Mo dongmei[18]discussed the necessity of thesaurus in the construction of medical ontology, and the construction flow of prevent medicine ontology, including ontology construction tool, method, and the convert from thesaurus to domain ontology, using racer inference machine to inference domain ontology.

Zhou xiaojia[19]proposed a solution of automatic extraction of temporal attributes of medical problems from Chinese narrative medical records based on conditional random

fields (CRF). Gao yonggang[20]used MI and its improved MI greedy optimization algorithm to achieve medical image feature selection; implement MI-based feature selection with support vector machine (SVMs) classification labels, so as to establish low-level visual features of medical images and high-level semantic features of the mapping between; solve medical image annotation in the "semantic gap" problem and all the characteristics of large amount of calculation involved in classification problems. Mi Yang[21]'s research is based on top ontology' integration of medical semantic annotation, it is of ontology's construction, ontology mapping and integration, semantic annotation, etc., explore the medical field relational structure of the top ontology semantic ontology semantic annotation working mechanism, try to establish a top-level ontology integration based on the medical domain ontology semantic annotation system model, empirical semantic annotation new ideas, explore ontology integration effect for semantic annotation role in promoting and strive to achieve from reality, to study and solve practical application of existing ontology semantic barriers, thereby improving medical information resources level semantic annotation. Tsatsaronis, George[22]proposed a method based on a Maximum Entropy approach, for annotating biomedical literature documents with terms from the Medical Subject Headings (MeSH).The experimental evaluation shows that the suggested Maximum Entropy approach for annotating biomedical documents with MeSH terms is highly accurate, robust to the ambiguity of terms, and can provide very good performance even when a very small number of training documents is used. Ruzicka, M; Svatek, V[23] proposed an annotation system—stepper, the system can extract information from literatures that are unique, authority, and rich in ambiguously expressed knowledge. Yu jiang de[24]proposed a method based on Conditional Random semantic role labeling, the method is based on shallow parsing, phrases or named entities is marked as the basic unit of the semantic annotation. CRFs model is used to predicate sentence semantic role labeling, experimental results show that: based on Conditional Random approach than a method based on maximum entropy better performance of the method in the semantic role labeling task gained 80.43% and 63.55% accuracy rate of recall. Volk, M[25]present a framework for concept-based cross-language information retrieval in the medical domain, which is under development in the MUCHMORE project.

Current ontology-based semantic annotation tools all have their own characteristics and scope, but the existing tools of the following deficiencies: 1. No tools support the latest W3C ontology language OWL; 2.most tools do not support ontology vocabulary expansion; 3.majority tools do not support simultaneous open, browse multiple ontology and ontology annotation using multiple pages; 4.mark all the tools only in English, does not support multi-language; 5.label objects tools is static content-based; 6. most of the tools used first create the content, after marked "two-step", only a few tools support content writing and semantic annotation simultaneously; 7. semantic annotation process of ontology query, assisted reasoning support, and the degree of automation of metadata generated is not enough [26].

**3. The theoretical basis of the Experiment.** To solve these problems, Liu Yao [27], etc.

According to the theory and bibliographic control conjugate theory using natural language processing (NLP) techniques and machine learning methods have been accepted domain knowledge, such as professional thesaurus, professional dictionaries, textbooks or professional the use of such authoritative works to reconstruct, build domain ontology, based on the development of Chinese literature resources semantic annotation technology for semantic annotation of the relevant literature, and in the relatively large number of semantic content, based on the combination of traditional organizational resources (thesauri, etc. ), through machine learning and other methods to generate the initial semantic metadata, and then with the help of the auxiliary platform library resources to achieve organizational processes and semantic semantic metadata system construction simultaneously, and semantic annotation semantic indexing literature and stored separately.

In this theory, the authors will take the skin diseases for example, where the Mesh terms in the auxiliary platform Ontology as the knowledge elements, its classification and property settings to be adjusted and modified in order to build the knowledge described in the prototype system. Prepared related literature, import the teaching materials into the system, to establish skin disease ontology automaticly, to verify and improve the method.
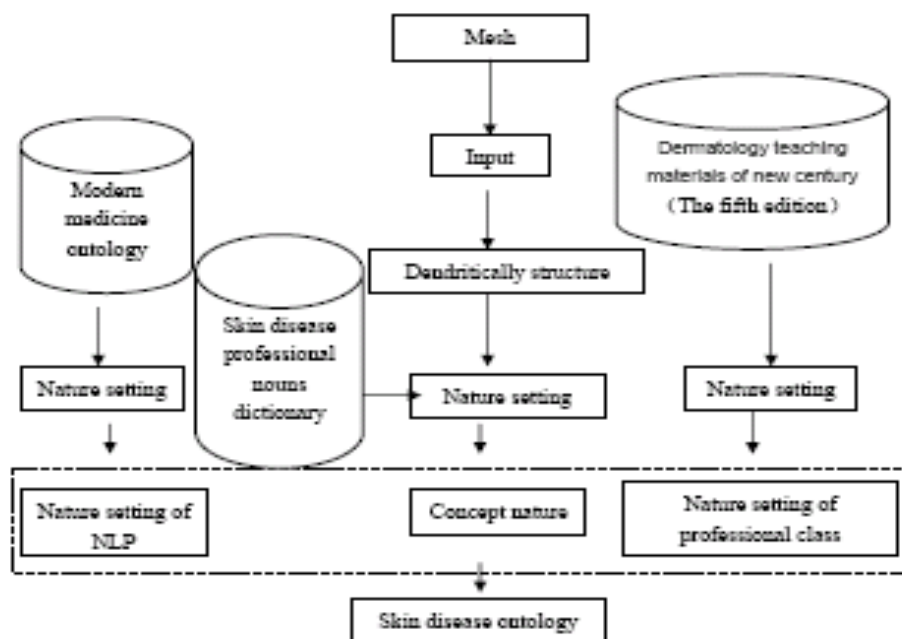


FIG.1.THE METHOD'S CONSTRUCTION SCHEDULE

## 4. The experiment's key technologie.

4.1. **Input of lexicon.** The lexicon we selected was Medical Subject Headings, which is editioned by NLM , we selected its part of skin disease and part that has relationships with the disease. This lexicon is structured, covers almost all skin diseases. We download the lexicon; sort it to unique format which the system asked for. The format of their hierarchical relationships, up and down relationship is represented with the Tab key

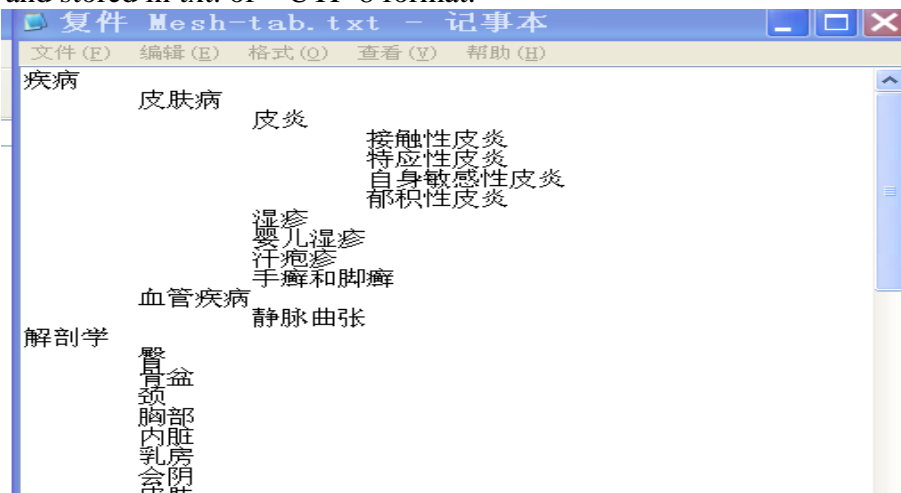hierarchy, and stored in txt. of  UTF-8 format.



FIG.2.LEXCION THAT THE SYSTEM REQUIRES



FIG.3.INPUT THE LEXICON TO THE SYSTEM

Input this lexicon to the platform， and after inputting the lexicon, the platform generates the class of dendritically structure as follows：
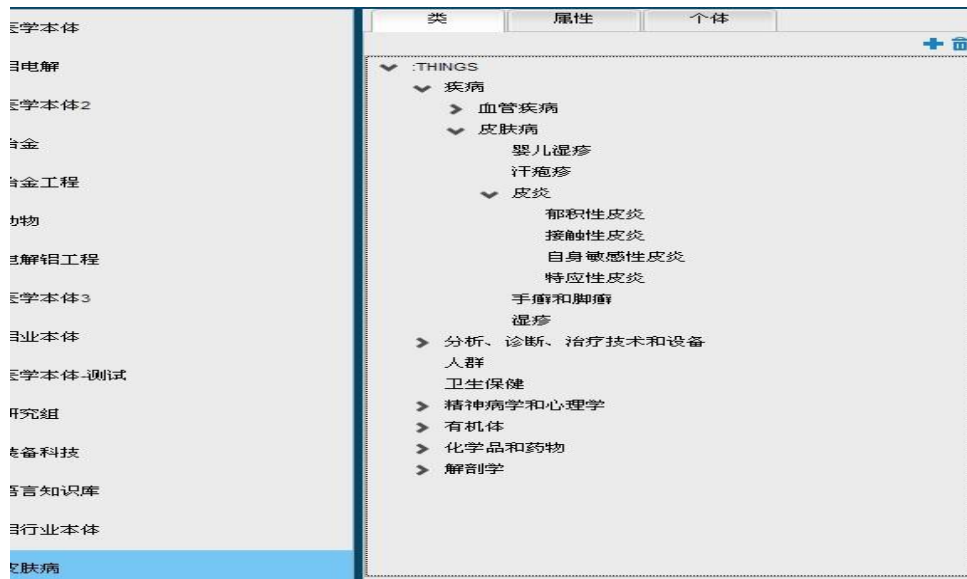
FIG.4. DENDRITICALLY STRUCTURE OF THE PLATFORM

4.2. **Nature setting.** Nature setting contains three parts .Part one is nature setting of the class, every class has public nature ,also has its unique class, for example disease has its public class FF，PA，X，XR, see related, code, also has its nature of Etiology and pathogenesis、clinical manifestations、diagnosis、treatment. So as the other class. after all classes 's nature has been seted , Part two is the setiing of the nature's domain of definition, every nature should be setted , Part three is the setting of the nature's actuating range, all the nature and nature with the number 2 should be setted.
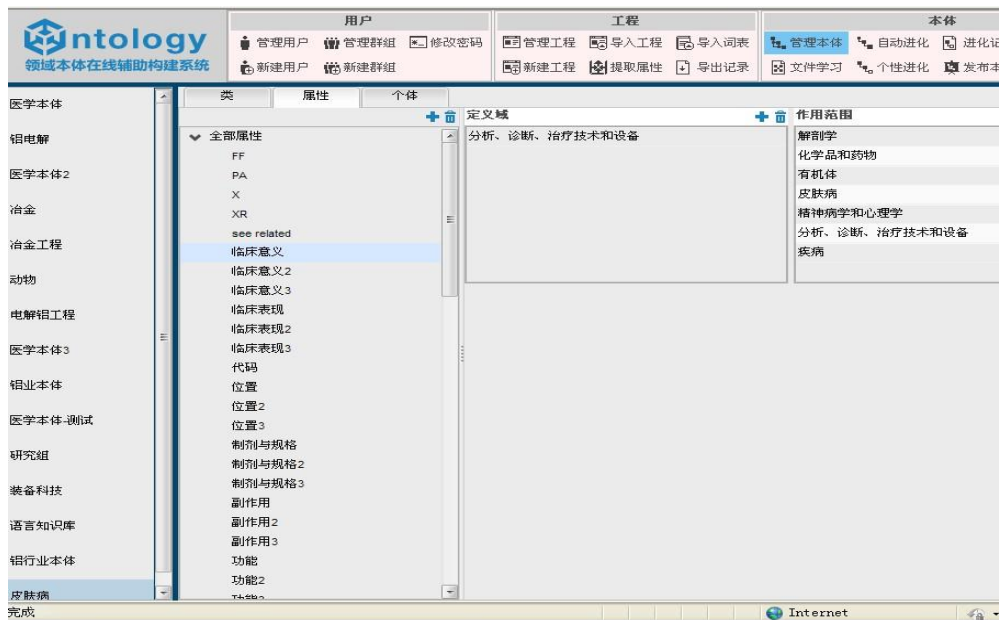


FIG.5.NATURE SETTING OF THE CLASS

23

4.3. **Format settings.** After the setting of the nature, to input the studying material, the studying material format should be setted to the format that the system requires, then the platform can identifies ,so that can use the platform to study the teaching materials.
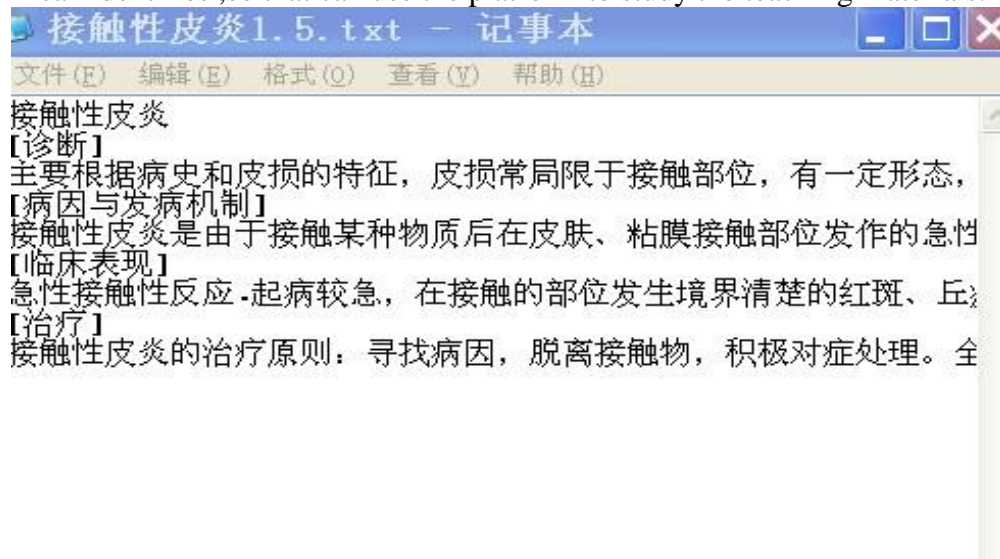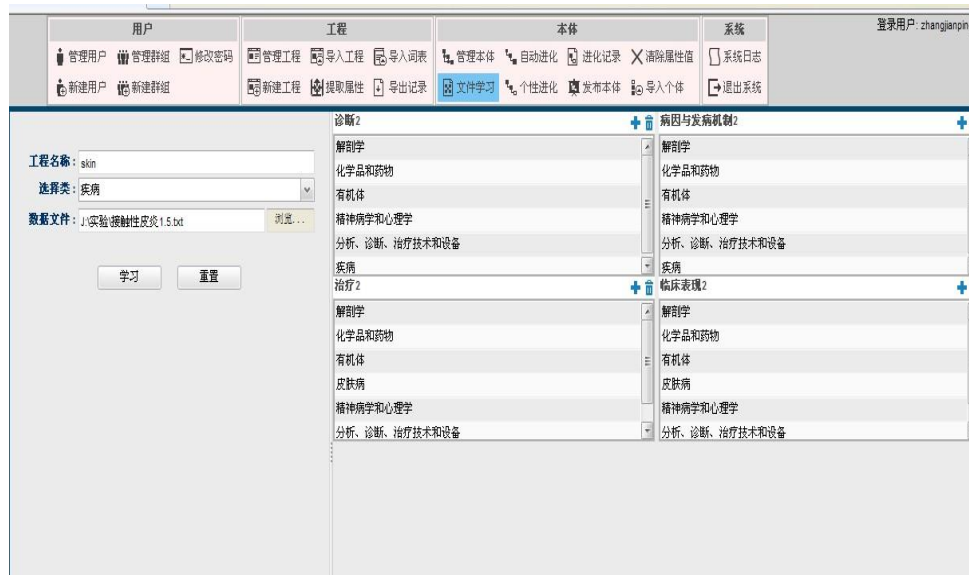


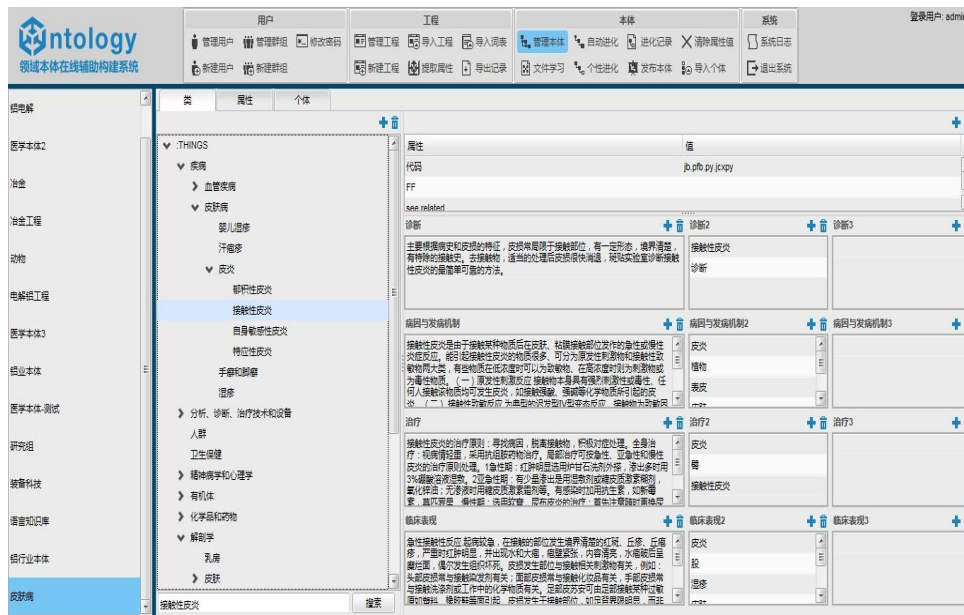FIG.6.STUDING MATERIAL' FORMAT



FIG.7.STUDYING METHOD

FIG.8.STUDYING RESULT

## 5. Conclusions.

The above is the studying of one material , it also can study other material , to Realize the function. This is about skin disease ontology automated build and annotation system development, this platform also need to further improve the indexing efficiency and accuracy.

**REFERENCES**

[1]   HAN Xianpei, ZHAO Jun.   Semantic Metadata Generation: A Method Based on Wikipedia［J］. Journal of Chinese Information Processing，2009，23 ( 2 ) : 108-114.

[2]   LIU Yao , SUI Zhi fang , HU Yong wei, et al.  Automatic construction of domain ontology [J]. Journal of Beijing University of Posts and Telecommunications, 2006, 29 ( Z2) : 65- 69.

[3]   http://semanticweb.org/wiki/Tools

[4]   Zhang Xiaolin. Semantic Web and Semantic-based Networked Information Retrieval, Journal of The China Society For Scientific and Technical Information ,2002, 14 (21).

[5]   HANDSCHUH S，et al. S-CREAM-Semi-automatic CREAtion of metadat［C］  13th International Conference， EKAW 2002 Sigüenza， Spain， 2002: 358-372.

[6]   CIMIANO P，et al.  Towards the self annotating Web[C] FELDMAN S I， et al.   Proceedings of WWW,2004:462-471.

[7]   CHIE CHUNG HUANG，et al. Using a Web based categorization approach to generate thematic

metadata from texts[J]. ACM Transactions on Asian Language Information Processing，2004，3（3）: 190-212.

[8] YANG H-C，LEE C-H. Automatic metadata generation for Web pages using a text mining approach ［C］International Workshop on Challenges in Web Information Retrieva1 and Integration，2005: 186-194.

[9] DINGLI A， et al. Automatic semantic annotation using unsupervised information extract on and integration ［C］GENNARIJ, et al. Proceedings K-CAP，2003.

[10] GRAUBITZ H, et al. Semantic tagging of domain-specific text documents with diasdem ［C］SAAKE G， et al. Proceedings of DBFusion 2001. USA: ACM，2001: 61-72.

[11] LI J，et al. Learning to generate semantic annotation for domain specific sentences ［C］GIL Y，et a1. Proceedings of KCAP，2001: 44-57.

[12] DILL S，et al. A case for automated largescale semantic annotation ［J］. Web Semantics: Science, Services and Agents onthe World Wide Web，2003，1（1）: 115-132.

[13] Liu Yao, Sui Zhifang, Zhou Yang .et al.Research on Automatic Construction of Chinese Traditional Medicine Ontology Concept' s Description Architecture. New Technology of Library and Information Service . 2008:5

[14] Meersman, R Creating ontologies for content representation - The OntoSeed suite

[15] Deng zi ping. Research and Development of a Ontology Automatic Generation System Oriented Medical Diagnosis. GuangDong University Of Technology[D]

[16] Fang an et al.Method Research of Constructing Clinical Disease Domain Ontology——A Case Study of Hand Foot and Mouth Disease Ontology. JOURNAL OF INTELLIGENCE.2009.v.28.No.11

[17] Fu qiang。 Research on Medical Domain Ontology Building Based on Thesaurus. Jilin University.[D]

[18] Mu Dongmei ,Fan Yi .Building and Reasoning of Digital Library's Domain Ontology:Taking Medicine Domain Ontology as an Example. Library and Information Service.[J]. Vol.51,No.8,August,2007

[19] ZHOU Xiao-Jia LI Hao-Min DUAN Hui-Long .et al.The Automatic Extraction of Temporal Relation from Chinese Narrative Medical Records Using Conditional Random Fields. Chinese Journal of Biomedical Engineering[J]. Vol. 29 No. 5, 2010

[20] Gao yonggang.Research and Application on Semantic Annotation Technologies of Medical Images[D]. Northwestern University. 2009

[21] Mi yang. Research on Semantic Annotation Based on Upper-level Ontology Integration in Medical Field.[D].JiLin university

[22] Tsatsaronis, George. A Maximum-Entropy approach for accurate document annotation in the biomedical domain.

[23] Ruzicka, M; Svatek, V .Stepper: Annotation and interactive stepwise transformation for knowledge-rich documents[C]. 14th International Conference on Knowledge Engineering and Knowledge Management

[24] Yu Jiangde,Fan Xiaozhong1 Pang Wenbo1. Semantic role labeling based on conditional random fields. Journal of Southeast University(English Edition.[J]

[25] Volk, M. Semantic annotation for concept-based cross-language medical information retrieval

[26] TAO Wan, LI Ping, LIAO Shu-mei. Analysis and summary of current ontology-based semantic annotation tools Journal of Anhui University of Technology and Science, 2005 ,20(2)

[27] Liu Yao et al., Research on the Method of Library Resources Organization Semantization Based on Content-format Interaction[J]. Information Studies:Theory & Application,. 2010 ,33(10)